
Stochastic L-BFGS Revisited: Improved Convergence Rates and Practical Acceleration Strategies

Renbo Zhao

Dept. ECE, ISEM & Math
National University of Singapore
elezren@nus.edu.sg

William B. Haskell

Dept. ISEM
National University of Singapore
isehwb@nus.edu.sg

Vincent Y. F. Tan

Dept. ECE & Math
National University of Singapore
vtan@nus.edu.sg

Abstract

We revisit the stochastic limited-memory BFGS (L-BFGS) algorithm. By proposing a new framework for analyzing convergence, we theoretically improve the (linear) convergence rates and computational complexities of the stochastic L-BFGS algorithms in previous works. In addition, we propose several practical acceleration strategies to speed up the empirical performance of such algorithms. We also provide theoretical analyses for most of the strategies. Experiments on large-scale logistic and ridge regression problems demonstrate that our proposed strategies yield significant improvements via-à-vis competing state-of-the-art algorithms.

1 INTRODUCTION

In this work, we are interested in the following (unconstrained) convex finite-sum minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left[f(\mathbf{x}) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \right], \quad (1)$$

where d and n denote the dimension of decision vector and the number of component functions respectively. Problems in the form of (1) play important roles in machine learning and signal processing. One important class of such problems is the *empirical risk minimization* (ERM) problem, where each f_i assumes the form

$$f_i(\mathbf{x}) \triangleq \ell(\mathbf{a}_i^T \mathbf{x}, b_i) + \lambda R(\mathbf{x}). \quad (2)$$

In (2), $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ denotes a smooth loss function, $R : \mathbb{R}^d \rightarrow \mathbb{R}_+$ a smooth convex regularizer (e.g., Tikhonov), $\lambda \geq 0$ the regularization weight and $\{(\mathbf{a}_i, b_i)\}_{i=1}^n \subseteq \mathbb{R}^{d+1}$ the set of feature-response pairs. Depending on the form of ℓ and R , many important machine learning problems are special cases of ERM, such as logistic regression, ridge regression and soft-margin support vector machine.

We focus on the case where both n and d are large, and f is ill-conditioned (i.e., the condition number of f is large).¹ In the context of ERM, this means the dataset $\{(\mathbf{a}_i, b_i)\}_{i=1}^n$ that defines (1) has large size and the feature vectors \mathbf{a}_i have high ambient dimension. However, the points typically belong to a low-dimensional manifold. Such a setting is particularly relevant in the big-data era, due to unprecedented data acquisition abilities.

When n is large, the computational cost incurred by the batch optimization methods (both first- and second-order) are prohibitive. Therefore, stochastic (randomized) optimization methods have become very popular in this setting. At each iteration, only a subset of component functions, rather than all of them, are processed. In this way, much more progress can be made towards global optima in the time that only allows batch methods to take a single step. When d is large, Newton- or quasi-Newton-based methods (both batch and stochastic) incur both high computational and storage complexities. Consequently only first-order and limited-memory quasi-newton methods (e.g., L-BFGS (Liu and Nocedal, 1989)) are feasible in this setting.

When both n and d are large, as in our setting, most of research efforts have been devoted to *stochastic first-order methods*, which include stochastic gradient descent (SGD) (Bottou, 1998; Bottou and LeCun, 2004) and its variance-reduced modifications (Defazio et al., 2014; Harikandeh et al., 2015; Johnson and Zhang, 2013; Schmidt et al., 2017). However, these methods do not make use of the curvature information. This limits their abilities to find highly accurate solutions for ill-conditioned problems. In order to incorporate the curvature information in the limited-memory setting, recently much progress have been made toward developing *stochastic L-BFGS* algorithm. A partial list of such works includes Bordes et al. (2009); Byrd et al. (2016); Gower et al. (2016); Mokhtari and Ribeiro (2015); Moritz et al. (2016); Schraudolph et al. (2007); Sohl-Dickstein et al. (2014). Among them, the SQN method (Byrd et al., 2016) innovatively makes use of

¹In this work, the condition number of a (strongly) convex function refers to that of its Hessian.

the subsampled Hessian-vector products to form the correction pairs (as opposed to using difference of stochastic gradients) and achieves better results than previous methods. However, the convergence rate is sublinear (in the strongly-convex case), similar to that of SGD. Later, Moritz et al. (2016) combines this method with *stochastic variance-reduced gradient* (SVRG) and proves linear convergence of the resulting algorithm. Gower et al. (2016) maintains the structure of this algorithm but incorporates the *block BFGS update* to collect more curvature information in the optimization process. Although the convergence rate of this new method remains similar to that in Moritz et al. (2016), experimental results have demonstrated practical speedup introduced by the block BFGS update.

1.1 MOTIVATIONS AND CONTRIBUTIONS

Our work can be motivated from both theory and practice. In terms of theory, although linear convergence has been shown for both algorithms in Moritz et al. (2016) and Gower et al. (2016), the convergence rates (and hence computational complexities) therein can be potentially further improved.² In terms of practice, in addition to block BFGS update, there may exist several other practical strategies that can potentially further accelerate³ the algorithm in Moritz et al. (2016). Based on these two aspects, our work makes the following contributions.

1) We propose a *coordinate transformation framework* to analyze the stochastic L-BFGS-type algorithms in Moritz et al. (2016) and Gower et al. (2016). Our framework yields a much improved convergence rate and computational complexity. Indeed, our approach can be applied to many other *stochastic second-order algorithms* as well (e.g., SQN algorithm in Byrd et al. (2016)); thus it has great generalizability and wide potential impacts. The essential idea of our method is to unify the analysis of stochastic first-order and second-order method, as a result, it *opens new avenues of designing and analyzing* other variants of stochastic second-order algorithms based on their first-order counterparts. Among these variants, two important examples include proximal (quasi-)Newton methods and momentum-based accelerated (quasi-)Newton methods.

2) We derive *two novel technical lemmas* on variance bound of stochastic gradient with *nonuniform mini-batch sampling* and spectral bound of the sequence of metric matrices $\{\mathbf{H}_r\}_{r \geq 0}$ in stochastic L-BFGS algorithms. These two lemmas are of independent interest. For example, the lemma on the spectral bound is equally applicable to batch L-BFGS algorithms.

²The analysis method in Gower et al. (2016) mainly follows that in Moritz et al. (2016), so we treat the analyses in both works in a unified manner.

³In this work, we refer “acceleration” to general strategies that speed up the algorithm, not necessarily the ones based on momentum methods.

3) We conduct a *systematic computational complexity analysis* for the stochastic L-BFGS algorithms. To the best of our knowledge, such an analysis is among the first for stochastic L-BFGS methods. Similar to our convergence analysis framework, our complexity analysis method can be easily generalized to other stochastic (quasi-)Newton methods as well.

4) We propose *several practical acceleration strategies* to speed up the convergence of the stochastic L-BFGS algorithm in Moritz et al. (2016). We demonstrate the efficacy of our strategies through numerical experiments on logistic and ridge regression problems. Theoretically, we also prove linear convergence for most strategies.

2 PRELIMINARIES

Notations. We use lowercase, bold lowercase and bold uppercase letters to denote scalars, vectors and matrices respectively. For a matrix $\mathbf{U} \in \mathbb{R}^{m_1 \times m_0}$, we denote its (p, q) -th entry as u_{pq} . For a function $f : \mathbb{R}^{m_1} \rightarrow \mathbb{R}^{m_2}$, define a function $f \circ \mathbf{U}$ such that $(f \circ \mathbf{U})(\mathbf{z}) \triangleq f(\mathbf{U}\mathbf{z})$, for any $\mathbf{z} \in \mathbb{R}^{m_0}$. We use \mathbb{N} to denote the set of natural numbers excluding zero. For any $n \in \mathbb{N}$, we define $[n] \triangleq \{1, \dots, n\}$ and $(n) \triangleq \{0, 1, \dots, n\}$. Accordingly, define $\mathcal{A}_{(n)} \triangleq \{\mathcal{A}_0, \dots, \mathcal{A}_n\}$. We use $\|\cdot\|$ to denote both the ℓ_2 norm of a vector and the spectral norm of a matrix. We use \mathbf{B} and \mathbf{H} (with subscripts and superscripts) to denote the approximate Hessian and approximate inverse Hessian in L-BFGS algorithms respectively, following the convention in Nocedal and Wright (2006). \mathbf{H} is also known as the *metric matrix* (Goldfarb, 1970). Sections and lemmas with indices beginning with ‘S’ will appear in the supplemental material.

Assumptions. We make two blanket assumptions on the component functions $\{f_i\}_{i \in [n]}$.

Assumption 1. For each $i \in [n]$, f_i is convex and twice differentiable on \mathbb{R}^d . For ERM problems (2), we assume these two properties are satisfied by the loss function ℓ in its first argument on \mathbb{R} and by the regularizer R on \mathbb{R}^d .

Assumption 2. For each $i \in [n]$, f_i is μ_i -strongly convex and L_i -smooth on \mathbb{R}^d , where $0 < \mu_i \leq L_i$.

Remark 1. Assumptions 1 and 2 are standard in the analysis of second-order optimization methods, for both deterministic and stochastic cases. The strong convexity of f_i in Assumption 2 ensures *positive* curvature at any point in \mathbb{R}^d , which in turn guarantees the well-definedness of the BFGS update. As a common practice in the literature (Byrd et al., 2016; Moritz et al., 2016), this condition can typically be enforced by adding a strongly convex regularizer (e.g., Tikhonov) to f_i . Due to the strong convexity, (1) has a unique solution, denoted as \mathbf{x}^* .

3 ALGORITHM

We present the pseudo-code of the algorithm in Moritz et al. (2016) in Algorithm 1. (Our acceleration strategies for this algorithm will be shown and discussed in Section 6.) In this algorithm, s denotes the outer iteration index, t the inner iteration index and r the index of metric matrices $\{\mathbf{H}_r\}_{r \geq 0}$. We use $\mathbf{x}_{s,t}$ and \mathbf{x}^s to denote an inner iterate and outer iterate respectively. Each outer iteration consists of m inner iterations. The only modification that we make on the original algorithm in Moritz et al. (2016) is sampling the index set $\mathcal{B}_{s,t} \subseteq [n]$ of size b with replacement *nonuniformly*. Specifically, the elements in $\mathcal{B}_{s,t}$ are sampled i.i.d. from a discrete distribution $P \triangleq (p_1, \dots, p_n)$, such that for any $i \in [n]$, $p_i = L_i / \sum_{i=1}^n L_i$. As will be seen in Lemma 3, compared to uniform sampling, nonuniform sampling leads to a better variance bound on the stochastic gradient $\mathbf{v}_{s,t}$. Using $\mathcal{B}_{s,t}$ and $\nabla f(\mathbf{x}^s)$, we then compute $\mathbf{v}_{s,t}$ according to (5) in Algorithm 1, where for any $\mathbf{x} \in \mathbb{R}^d$,

$$\nabla f_{\mathcal{B}_{s,t}}(\mathbf{x}) \triangleq \frac{1}{b} \sum_{i \in \mathcal{B}_{s,t}} \frac{1}{np_i} \nabla f_i(\mathbf{x}). \quad (3)$$

Next we compute the search direction formed by the product of the metric matrix \mathbf{H}_r and $\mathbf{v}_{s,t}$. In the limited-memory setting, we do not store \mathbf{H}_r , but only a set \mathcal{H}_r , which consists of $M' \triangleq \min\{r, M\}$ recent *correction pairs* $\{(\mathbf{s}_j, \mathbf{y}_j)\}_{j=r-M'+1}^r$. The correction pairs are constructed from the averaged past iterates $\{\bar{\mathbf{x}}_r\}_{r \geq 0}$ (defined in line 13), and thus contains the ‘‘smoothed’’ curvature information in recent iterates. Following Byrd et al. (2016), in computing \mathbf{y}_r , we first sample an index set $\mathcal{T}_r \subseteq [n]$ of size b_H uniformly without replacement, and then let $\mathbf{y}_r = \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \mathbf{s}_r$, where

$$\nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \triangleq \frac{1}{b_H} \sum_{i \in \mathcal{T}_r} \nabla^2 f_i(\bar{\mathbf{x}}_r), \quad (4)$$

denotes the sub-sampled Hessian at $\bar{\mathbf{x}}_r$.

In computing $\mathbf{H}_r \mathbf{v}_{s,t}$, a direct approach would be to compute \mathbf{H}_r first and form the product. Computing \mathbf{H}_r involves applying M' BFGS updates to the matrix $\mathbf{H}_r^{(r-M')} \triangleq (\mathbf{s}_r^T \mathbf{y}_r) / \|\mathbf{y}_r\|^2 \mathbf{I}$ using $\{(\mathbf{s}_j, \mathbf{y}_j)\}_{j=r-M'+1}^r$. For each $k \in \{r-M'+1, \dots, r\}$, the update is

$$\mathbf{H}_r^{(k)} = \left(\mathbf{I} - \frac{\mathbf{y}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) \mathbf{H}_r^{(k-1)} \left(\mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}.$$

Finally we set $\mathbf{H}_r = \mathbf{H}_r^{(r)}$. Instead of using this direct approach, we adopt the two-loop recursion algorithm (Nocedal and Wright, 2006, Algorithm 7.4) to compute $\mathbf{H}_r \mathbf{v}_{s,t}$ as a whole. This approach has the same effect as the direct one, but with much reduced computation.

At the end of each outer iteration, the starting point of next outer iteration, \mathbf{x}^{s+1} is either uniformly sampled (option

Algorithm 1 Stochastic L-BFGS Algorithm with Nonuniform Mini-batch Sampling

- 1: **Input:** Initial decision vector \mathbf{x}^0 , mini-batch sizes b and b_H , parameters m, M and Υ , step-size η , termination threshold ϵ
 - 2: **Initialize** $s := 0, r := 0, \bar{\mathbf{x}}_0 = \mathbf{0}, \mathbf{H}_0 := \mathbf{I}, \mathcal{H}_0 := \emptyset$
 - 3: **Repeat**
 - 4: Compute a full gradient $\mathbf{g}_s \triangleq \nabla f(\mathbf{x}^s)$
 - 5: $\mathbf{x}_{s,0} := \mathbf{x}^s$
 - 6: **for** $t = 0, 1, \dots, m-1$
 - 7: Sample a set $\mathcal{B}_{s,t}$ with size b
 - 8: Compute a variance-reduced gradient

$$\mathbf{v}_{s,t} := \nabla f_{\mathcal{B}_{s,t}}(\mathbf{x}_{s,t}) - \nabla f_{\mathcal{B}_{s,t}}(\mathbf{x}^s) + \mathbf{g}_s \quad (5)$$
 - 9: Compute $\mathbf{H}_r \mathbf{v}_{s,t}$ from \mathcal{H}_r and $\mathbf{v}_{s,t}$
 - 10: $\mathbf{x}_{s,t+1} := \mathbf{x}_{s,t} - \eta \mathbf{H}_r \mathbf{v}_{s,t}$
 - 11: **if** $sm + t > 0$ **and** $(sm + t) \equiv 0 \pmod{\Upsilon}$
 - 12: $r := r + 1$
 - 13: $\bar{\mathbf{x}}_r := \frac{1}{\Upsilon} \sum_{l=t-\Upsilon+1}^t \mathbf{x}_{s,l}$
 - 14: Sample a set \mathcal{T}_r with size b_H
 - 15: $\mathbf{s}_r := \bar{\mathbf{x}}_r - \bar{\mathbf{x}}_{r-1}, \mathbf{y}_r := \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \mathbf{s}_r$
 - 16: Update $\mathcal{H}_r := \{(\mathbf{s}_j, \mathbf{y}_j)\}_{j=r-M'+1}^r$
 - 17: **end if**
 - 18: **end for**
 - 19: **Option I:** Sample τ_s uniformly randomly from $[m]$ and set $\mathbf{x}^{s+1} := \mathbf{x}_{s,\tau_s}$
 - 20: **Option II:** $\mathbf{x}^{s+1} := \frac{1}{m} \sum_{t=1}^m \mathbf{x}_{s,t}$
 - 21: $s := s + 1$
 - 22: **Until** $|f(\mathbf{x}^s) - f(\mathbf{x}^{s-1})| < \epsilon$
 - 23: **Output:** \mathbf{x}^s
-

I) or averaged (option II) from all the past inner iterates $\{\mathbf{x}_{s,t}\}_{t \in [m]}$. As shown in Theorem 1, these two options can be analyzed in a unified manner.

4 CONVERGENCE ANALYSIS

4.1 DEFINITIONS

Let $\{\underline{i}_j\}_{j \in [n]}$ and $\{\bar{i}_j\}_{j \in [n]}$ be permutations of $[n]$ such that $\mu_{\min} \triangleq \mu_{\underline{i}_1} \leq \dots \leq \mu_{\underline{i}_n}$ and $L_{\bar{i}_1} \leq \dots \leq L_{\bar{i}_n} \triangleq L_{\max}$. Given any $\tilde{n} \in [n]$, define

$$\bar{\mu}_{\tilde{n}} \triangleq \frac{1}{\tilde{n}} \sum_{j=1}^{\tilde{n}} \mu_{\underline{i}_j} \quad \text{and} \quad \bar{L}_{\tilde{n}} \triangleq \frac{1}{\tilde{n}} \sum_{j=n-\tilde{n}+1}^n L_{\bar{i}_j}. \quad (6)$$

Accordingly, define $\kappa_{\max} \triangleq L_{\max} / \mu_{\min}$ and $\kappa_{\tilde{n}} \triangleq \bar{L}_{\tilde{n}} / \bar{\mu}_{\tilde{n}}$. In particular, define $\bar{\mu} \triangleq \bar{\mu}_{\tilde{n}}, \bar{L} \triangleq \bar{L}_{\tilde{n}}$ and $\kappa \triangleq \bar{L} / \bar{\mu}$. We also define a filtration $\{\mathcal{F}_{s,t}\}_{s \geq 0, t \in (m-1]}$ such that $\mathcal{F}_{s,t}$ contains all the information up to the time (s, t) . Formally, $\mathcal{F}_{s,t} \triangleq \sigma(\{\tau_j\}_{j=0}^{s-1} \cup \{\mathcal{B}_{i,j}\}_{i \in (s-1), j \in (m-1)} \cup \{\mathcal{B}_{s,j}\}_{j=0}^{t-1} \cup \{\mathcal{T}_j\}_{j=0}^{\lfloor (sm+t)/L \rfloor})$, where $\sigma(\{x_j\}_{j=1}^n)$ denotes

the σ -algebra generated by random variables $\{x_j\}_{j=1}^n$. We also define $\mathcal{F}_s \triangleq \mathcal{F}_{s,0}$.

To introduce our coordinate transformation framework, we define some transforms of variables appearing in Algorithm 1. Specifically, for any $s, t, r \geq 0$, define $\tilde{\mathbf{x}}_{s,t,r} \triangleq \mathbf{H}_r^{-1/2} \mathbf{x}_{s,t}$, $\tilde{\mathbf{x}}^{s,r} \triangleq \mathbf{H}_r^{-1/2} \mathbf{x}^s$, $\tilde{\mathbf{x}}_r^* \triangleq \mathbf{H}_r^{-1/2} \mathbf{x}^*$ and $\tilde{\mathbf{v}}_{s,t,r} \triangleq \mathbf{H}_r^{1/2} \mathbf{v}_{s,t}$.⁴ We also define transformed functions $\tilde{f}_{i,r} \triangleq f_i \circ \mathbf{H}_r^{1/2}$ and $\tilde{f}_r \triangleq \frac{1}{n} \sum_{i=1}^n \tilde{f}_{i,r}$, for any $i \in [n]$ and $r \geq 0$.

4.2 PRELIMINARY LEMMAS

From the definitions of transformed variables and functions in Section 4.1, we immediately have the following lemmas.

Lemma 1. *For any $s, t, r \geq 0$ and $i \in [n]$, we have*

$$\tilde{f}_{i,r}(\tilde{\mathbf{x}}_{s,t,r}) = f_i(\mathbf{x}_{s,t}), \quad (7)$$

$$\nabla \tilde{f}_{i,r}(\tilde{\mathbf{x}}_{s,t,r}) = \mathbf{H}_r^{1/2} \nabla f_i(\mathbf{x}_{s,t}), \quad (8)$$

$$\nabla^2 \tilde{f}_{i,r}(\tilde{\mathbf{x}}_{s,t,r}) = \mathbf{H}_r^{1/2} \nabla^2 f_i(\mathbf{x}_{s,t}) \mathbf{H}_r^{1/2}. \quad (9)$$

Lemma 2. *If there exist $0 < \gamma' \leq \Gamma'$ such that $\gamma' \mathbf{I} \preceq \mathbf{H}_r \preceq \Gamma' \mathbf{I}$ for all $r \geq 0$, then for any $i \in [n]$ and $r \geq 0$, $\tilde{f}_{i,r}$ is twice differentiable, $(\mu_i \gamma')$ -strongly convex and $(L_i \Gamma')$ -smooth on \mathbb{R}^d . Consequently, \tilde{f}_r is twice differentiable, $(\gamma' \bar{\mu})$ -strongly convex and $(\Gamma' \bar{L})$ -smooth on \mathbb{R}^d .*

Next we derive two other lemmas that will not only be used in the analysis later, but have potential to be applied to more general problem settings. Specifically, Lemma 3 can be applied to any stochastic optimization algorithms based on SVRG and Lemma 4 can be applied to any finite-sum minimization algorithms based on L-BFGS methods (not necessarily stochastic in nature). The proofs for Lemmas 3 and 4 are deferred to Sections S-3 and S-4 respectively.

Lemma 3 (Variance bound of $\mathbf{v}_{s,t}$). *In Algorithm 1, we have $\mathbb{E}_{\mathcal{B}_{s,t}} [\mathbf{v}_{s,t} | \mathcal{F}_{s,t}] = \nabla f(\mathbf{x}_{s,t})$ and*

$$\begin{aligned} & \mathbb{E}_{\mathcal{B}_{s,t}} \left[\|\mathbf{v}_{s,t} - \nabla f(\mathbf{x}_{s,t})\|^2 | \mathcal{F}_{s,t} \right] \\ & \leq \frac{4\bar{L}}{b} (f(\mathbf{x}_{s,t}) - f(\mathbf{x}^*) + f(\mathbf{x}^s) - f(\mathbf{x}^*)). \end{aligned} \quad (10)$$

Remark 2. In previous works (Gower et al., 2016; Moritz et al., 2016), a uniform mini-batch sampling of $\mathcal{B}_{s,t}$ was employed, and different variance bounds of $\mathbf{v}_{s,t}$ were derived. In Moritz et al. (2016, Lemma 6), the bound was

$$4L_{\max} (f(\mathbf{x}_{s,t}) - f(\mathbf{x}^*) + f(\mathbf{x}^s) - f(\mathbf{x}^*)). \quad (11)$$

In Gower et al. (2016, Lemma 2), this bound was slightly improved to be

$$4L_{\max} ((f(\mathbf{x}_{s,t}) - f(\mathbf{x}^*)) + (1 - 1/\kappa_{\max}) (f(\mathbf{x}^s) - f(\mathbf{x}^*))). \quad (12)$$

⁴As will be shown in Lemma 4, for any $r \geq 0$, $\mathbf{H}_r \succeq \gamma \mathbf{I}$ for some $\gamma > 0$. Therefore, $\mathbf{H}_r^{1/2}$ (and $\mathbf{H}_r^{-1/2}$) are well-defined.

However, both of these bounds fail to capture the dependence on the mini-batch size b . In contrast, in this work we consider a nonuniform mini-batch sampling (with replacement). Due to division by b and $\bar{L} \leq L_{\max}$ (indeed in many cases, $\bar{L} \ll L_{\max}$), our bound in (10) is much superior to (11) and (12). As will be seen in Theorem 1, our better bound (10) leads to a faster (linear) convergence rate of Algorithm 1.

Lemma 4 (Uniform Spectral Bound of $\{\mathbf{H}_r\}_{r \geq 0}$). *The spectra of $\{\mathbf{H}_r\}_{r \geq 0}$ are uniformly bounded, i.e., for each $r \geq 0$, $\gamma \mathbf{I} \preceq \mathbf{H}_r \preceq \Gamma \mathbf{I}$, where⁵*

$$\gamma \triangleq \frac{1}{(M+1)\bar{L}_{b_H}} \quad \text{and} \quad \Gamma \triangleq \frac{\kappa_{b_H}^{M+1}}{\bar{\mu}_{b_H}(\kappa_{b_H} - 1)}. \quad (13)$$

Remark 3. In Byrd et al. (2016) and Moritz et al. (2016), the authors make use of a classical technique in Liu and Nocedal (1989) to derive a different uniform spectral bound of $\{\mathbf{H}_r\}_{r \geq 0}$. Their technique involves applying $\text{tr}(\cdot)$ and $\det(\cdot)$ recursively to the BFGS update rule

$$\mathbf{B}_r^{(k)} = \mathbf{B}_r^{(k-1)} - \frac{\mathbf{B}_r^{(k-1)} \mathbf{s}_k \mathbf{s}_k^T \mathbf{B}_r^{(k-1)}}{\mathbf{s}_k^T \mathbf{B}_r^{(k-1)} \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k}, \quad (14)$$

where $\mathbf{B}_r^{(k)} \triangleq (\mathbf{H}_r^{(k)})^{-1}$ denotes the approximate Hessian matrix at step k in the reconstruction of $\mathbf{B}_r \triangleq (\mathbf{H}_r)^{-1}$. The lower and upper bounds derived by this technique are

$$\tilde{\gamma} = \frac{1}{(d+M)L_{\max}} \quad \text{and} \quad \tilde{\Gamma} = (d+M)^{d+M-1} \frac{\kappa_{\max}^{d+M-1}}{\mu_{\min}}$$

respectively. As will be seen in Proposition 1, the overall computational complexity of Algorithm 1 heavily depends on the estimated (uniform) condition number of $\{\mathbf{H}_r\}_{r \geq 0}$. Therefore, it is instructive to compute this quantity for both (γ, Γ) and $(\tilde{\gamma}, \tilde{\Gamma})$ as

$$\kappa_H \triangleq \frac{\Gamma}{\gamma} = (M+1) \frac{\kappa_{b_H}^{M+2}}{\kappa_{b_H} - 1} \approx (M+1) \kappa_{b_H}^{M+1}, \quad (15)$$

$$\tilde{\kappa}_H \triangleq \frac{\tilde{\Gamma}}{\tilde{\gamma}} = (M+d)^{M+d} \kappa_{\max}^{M+d}, \quad (16)$$

where the approximation in (15) follows from $\kappa_{b_H} \gg 1$ (see Footnote 5). By comparing (15) and (16), we notice our estimate for the condition number of $\{\mathbf{H}_r\}_{r \geq 0}$, namely κ_H , is smaller than those in Byrd et al. (2016) and Moritz et al. (2016), namely $\tilde{\kappa}_H$, in several aspects. First, κ_H does not grow (exponentially) with the data dimension d . Second, κ_H depends on κ_{b_H} , which is usually much smaller than κ_{\max} . Third, even if we set $d = 1$ in (16), the factor $M+1$ in (15) is much smaller than the factor $(M+1)^{M+1}$ in (16). As a result, our improved estimate of the condition number of $\{\mathbf{H}_r\}_{r \geq 0}$ will lead to a much better computational complexity estimate (see Proposition 1).

⁵We assume $\kappa_{b_H} > 1$ for any $b_H \in [n]$ since we focus on the setting where f is ill-conditioned, i.e., $\kappa_{b_H} \geq \kappa \gg 1$. If $\kappa_{b_H} = 1$ for some $b_H \in [n]$, then $\Gamma = (\kappa_{b_H}^M + M)/\bar{\mu}_{b_H}$ and γ remains the same. The proof for this case can be straightforwardly adapted from that in Section S-4.

4.3 MAIN RESULT

Theorem 1 presents our main convergence result. For the complete proof, we refer readers to Section S-2.

Theorem 1. *In Algorithm 1, if we choose $\eta < \min\{b/12, 1\}/(\Gamma\bar{L})$ and m sufficiently large, then with either option I or II, we have*

$$\mathbb{E}[f(\mathbf{x}^s) - f(\mathbf{x}^*)] \leq \rho^s (f(\mathbf{x}^0) - f(\mathbf{x}^*)), \text{ where} \quad (17)$$

$$\rho = \frac{b}{\gamma\bar{\mu}m\eta(b - 4\eta\Gamma\bar{L})} + \frac{4\eta\Gamma\bar{L}}{b - 4\eta\Gamma\bar{L}} \left(1 + \frac{1}{m}\right) < 1. \quad (18)$$

Remark 4. We compare our linear convergence rate ρ in (18) with those in Moritz et al. (2016) and Gower et al. (2016). Since the convergence rates in these two works are almost the same, we use the rate in Moritz et al. (2016) for comparison. The linear rate $\tilde{\rho}$ in Moritz et al. (2016) equals

$$\frac{1}{2\tilde{\gamma}\mu_{\min}m\eta(1 - \eta\tilde{\Gamma}L_{\max}\kappa_{\max}\tilde{\kappa}_H)} + \frac{\eta\tilde{\Gamma}L_{\max}\kappa_{\max}\tilde{\kappa}_H}{1 - \eta\tilde{\Gamma}L_{\max}\kappa_{\max}\tilde{\kappa}_H}.$$

For simplicity, if we let $b = 1$, $\mu_{\min} = \bar{\mu}$, $L_{\max} = \bar{L}$, $\tilde{\gamma} = \gamma$, $\tilde{\Gamma} = \Gamma$ and ignore other constant factors,⁶ we notice that there is an additional multiplicative factor κ_{κ_H} associated with $\eta\Gamma\bar{L}$ in $\tilde{\rho}$. As a result $\tilde{\rho} > \rho$. A more direct way to observe the detrimental effects of this additional κ_{κ_H} is to compare the computational complexities resulting from ρ and $\tilde{\rho}$. See Remark 6 for details. The reason that we manage to avoid this factor in our rate ρ is precisely because we adopt the *coordinate transformation framework* in our analysis (see proof sketch below). Specifically, by absorbing the sequence of metric matrices $\{\mathbf{H}_r\}_{r \geq 0}$ into decision vectors and functions, we are able to proceed through bounding the (expected squared Euclidean) distance between $\tilde{\mathbf{x}}_{s,t+1,r}$ and $\tilde{\mathbf{x}}_r^*$, instead of directly bounding $f(\mathbf{x}_{s,t+1})$ via the smoothness property of f (cf. proof of Theorem 7 in Moritz et al. (2016)). Thus in our analysis, we avoid the additional appearance of \bar{L} and Γ (which leads to the additional factor κ_{κ_H}).

Proof Sketch. Fix an outer iteration s and consider an inner iteration t . Define $r \triangleq \lfloor (sm + t)/L \rfloor$,⁷ then the iteration in (10) becomes

$$\tilde{\mathbf{x}}_{s,t+1,r} = \tilde{\mathbf{x}}_{s,t,r} - \eta\tilde{\mathbf{v}}_{s,t,r}. \quad (19)$$

From Lemmas 1, 3 and 4, we have $\mathbb{E}_{\mathcal{B}_{s,t}}[\tilde{\mathbf{v}}_{s,t,r} | \mathcal{F}_{s,t}] = \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r})$ and

$$\mathbb{E}_{\mathcal{B}_{s,t}} \left[\|\tilde{\mathbf{v}}_{s,t,r} - \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r})\|^2 | \mathcal{F}_{s,t} \right] \leq \frac{4\eta\bar{L}}{b} (\tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*)), \quad (20)$$

⁶However, bear in mind that by doing all these substitutions, $\tilde{\rho}$ has already been much improved.

⁷To avoid cluttered notations, we omit showing the dependence of r on s and t .

where $r' \triangleq \lfloor sm/L \rfloor$. Define $\tilde{\delta}_{s,t,r} \triangleq \tilde{\mathbf{v}}_{s,t,r} - \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r})$. Using $\eta \leq 1/(\Gamma\bar{L})$ and Lemma 2, we are able to show that

$$\|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 \leq \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 - 2\eta(\tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*)) + 2\eta^2\|\tilde{\delta}_{s,t,r}\|^2 - 2\eta \left\langle \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \right\rangle.$$

Taking expectation w.r.t. $\mathcal{B}_{s,t}$ and using (S-8), we have

$$\mathbb{E}_{\mathcal{B}_{s,t}} \left[\|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 + 2\eta(\tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*)) | \mathcal{F}_{s,t} \right] \leq \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 + \frac{8}{b}\Gamma\bar{L}\eta^2(\tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*)). \quad (21)$$

Telescoping (S-20) over $t \in (0, m-1]$ and noting that both options I and II in Algorithm 1 yield

$$\frac{1}{m} \sum_{t=1}^m \mathbb{E}_{\mathcal{B}_{s,(t-1)}} \left[\tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) | \mathcal{F}_{s,t-1} \right] \geq \mathbb{E}_{\mathcal{B}_{s,(m-1)}} \left[\tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) | \mathcal{F}_s \right], \quad (22)$$

where $r'' \triangleq \lfloor (s+1)m/L \rfloor$, we have

$$2m\eta \left(1 - \frac{4}{b}\Gamma\bar{L}\eta\right) \mathbb{E}_{\mathcal{B}_{s,(m-1)}} \left[\tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) | \mathcal{F}_s \right] \leq \|\tilde{\mathbf{x}}^{s,r'} - \tilde{\mathbf{x}}_{r'}^*\|^2 + \frac{8}{b}\Gamma\bar{L}\eta^2(1+m)(\tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*)).$$

Using $(\gamma\bar{\mu})$ -strong convexity of $\tilde{f}_{r'}$, we can derive the bound $\|\tilde{\mathbf{x}}^{s,r'} - \tilde{\mathbf{x}}_{r'}^*\|^2 \leq \frac{2}{\gamma\bar{\mu}}(\tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*))$. Finally, using Lemma 1 and rearranging, we have

$$\mathbb{E}[f(\mathbf{x}^{s+1}) - f(\mathbf{x}^*) | \mathcal{F}_s] \leq \rho(f(\mathbf{x}^s) - f(\mathbf{x}^*)). \quad (23)$$

Taking expectation on both sides, we reach (17). \square

5 COMPLEXITY ANALYSIS

In this section we provide a systematic framework for analyzing the computational complexity of the stochastic L-BFGS algorithm in Algorithm 1. Note that our framework can be easily generalized to other stochastic second-order algorithms, e.g., SQN algorithm in Byrd et al. (2016). To begin with, we make two additional assumptions.

Assumption 3. *For any $\mathbf{x} \in \mathbb{R}^d$ and $i \in [n]$, the gradient $\nabla f_i(\mathbf{x})$ can be computed in $O(d)$ operations.⁸*

Assumption 4. *For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $i \in [n]$, the Hessian-vector product $\nabla^2 f_i(\mathbf{x})\mathbf{y}$ can be computed in $O(d)$ operations.*

Remark 5. These two assumptions are naturally satisfied for ERM problems (2) with Tikhonov regularization. For these problems, $R(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$ and

$$\nabla \ell(\mathbf{a}_i^T \mathbf{x}, b_i) = \ell'(\mathbf{a}_i^T \mathbf{x}, b_i)\mathbf{a}_i + \lambda \mathbf{x}, \quad (24)$$

$$\nabla^2 \ell(\mathbf{a}_i^T \mathbf{x}, b_i)\mathbf{y} = \ell''(\mathbf{a}_i^T \mathbf{x}, b_i)(\mathbf{a}_i^T \mathbf{y})\mathbf{a}_i + \lambda \mathbf{y}, \quad (25)$$

⁸An operation refers to evaluation of an elementary function, such as addition, multiplication and logarithm.

where $\ell'(\cdot, \cdot)$ and $\ell''(\cdot, \cdot)$ are first and second derivatives of $\ell(\cdot, \cdot)$ w.r.t. the first argument. We easily see that the right-hand sides of both (24) and (25) can be computed in $O(d)$ operations.

From Algorithm 1, we observe that its total computational cost C can be split into three parts. The first part C_1 involves computing the variance-reduced gradient $\mathbf{v}_{s,t}$ in (5), the second part C_2 involves computing $\mathbf{H}_r \mathbf{v}_{s,t}$ (via two-loop recursion) in line 9, and the third part C_3 involves computing the correction pair $(\mathbf{s}_r, \mathbf{y}_r)$ in line 15.

Proposition 1. *Assume Assumptions 1 to 4. In Algorithm 1,*

$$C_1 = O((n + \kappa\kappa_H)d \log(1/\epsilon)), \quad (26)$$

$$C_2 = O(\kappa\kappa_H d \log(1/\epsilon)), \quad (27)$$

$$C_3 = O(d \log(1/\epsilon)). \quad (28)$$

Thus the total computational cost $C \triangleq \sum_{i=1}^3 C_i$ equals

$$O((n + \kappa\kappa_H) d \log(1/\epsilon)). \quad (29)$$

Proof. We leverage techniques that have become standard in the SVRG literature (e.g., Xiao and Zhang (2014)). In (18), if we choose $\eta = \theta b / (\Gamma \bar{L})$ for some $0 < \theta < 1/12$, $m = \theta' \kappa\kappa_H / b$ for some sufficiently large positive constant θ' , and use $1 + 1/m \leq 2$, then $\rho = \frac{1}{\theta' \theta (1-4\theta)} + \frac{8\theta}{1-4\theta} < 1$. As a result, the required number of outer iterations to achieve ϵ -suboptimality is $O(\log(1/\epsilon))$. Thus (26) follows from Assumption 3 and that $2mb$ gradients (of component functions) are computed in each inner iteration. If we further choose $M = \Theta(b)$, then (27) follows from the fact that two-loop recursion can be done in $O(Md)$ time (Nocedal and Wright, 2006, Chapter 7). Lastly, if we choose $b_H = \Theta(\Upsilon)$, then we obtain (28) using Assumption 4. \square

Remark 6. Following a similar argument, we can deduce the total complexity estimate \tilde{C} based on the linear rate $\tilde{\rho}$ (see Remark 4) derived in Moritz et al. (2016) as

$$\tilde{C} = O((n + b(\kappa_{\max} \tilde{\kappa}_H)^2) d \log(1/\epsilon)). \quad (30)$$

Compared with \tilde{C} , we observe that our complexity estimate C in (29) is much better, in several aspects. First, the dependence of C on the condition number $\kappa\kappa_H$ is linear, rather than quadratic. The quadratic dependence of $\kappa_{\max} \tilde{\kappa}_H$ in \tilde{C} is precisely caused by the additional $\kappa_{\max} \tilde{\kappa}_H$ in $\tilde{\rho}$ (see Remark 4). Second, C is independent of the min-batch size b . The appearance of b in \tilde{C} is a result of the loose bound on variance of $\mathbf{v}_{s,t}$ (cf. (11) and (12)). Third, the condition number $\kappa\kappa_H$ in C is much more smaller than $\kappa_{\max} \tilde{\kappa}_H$ in \tilde{C} for ill-conditioned problems. This is a result of non-uniform sampling of $\mathcal{B}_{s,t}$ and our improved bound on the spectra of $\{\mathbf{H}_r\}_{r \geq 0}$ (see Lemma 4).

Remark 7. As our coordinate transformation framework unifies the design and analysis of stochastic first- and second-order algorithms, we believe that momentum-based

acceleration techniques for stochastic first-order methods (Allen-Zhu, 2016; Lin et al., 2015) can be applied to Algorithm 1 as well. (Details are left to future work.) In this case, the dependence on $\kappa\kappa_H$ in C may be further improved to $\sqrt{\kappa\kappa_H}$ (Allen-Zhu, 2016).

6 ACCELERATION STRATEGIES

In this section, we propose three practical acceleration strategies. We follow the notations in Section 3 and Algorithm 1. As will be shown in Section 7.1, all of these strategies result in faster convergence in practice. For first and second strategies, we also provide their theoretical analyses in Propositions 2 and 3 respectively. See Sections S-5 and S-6 for the proofs of these two propositions.

6.1 GEOMETRIC SAMPLING AND AVERAGING

Instead of choosing \mathbf{x}^{s+1} according to option I or II in Algorithm 1, inspired by Konečný and Richtárik (2013), we can introduce a “forgetting” effect by considering two other sampling or averaging schemes:

option III: Sample τ_s randomly from $[m]$ according to the distribution $Q \triangleq (\beta^{m-1}/c, \beta^{m-2}/c, \dots, 1/c)$ and set $\mathbf{x}^{s+1} := \mathbf{x}_{s, \tau_s}$,

option IV: $\mathbf{x}^{s+1} := \frac{1}{c} \sum_{t=1}^m \beta^{m-t} \mathbf{x}_{s,t}$,

where $0 < \beta \leq 1 - \eta\gamma\bar{\mu} < 1$ and the normalization constant $c \triangleq \sum_{t=1}^m \beta^{m-t}$. Since $\beta \in (0, 1)$, we observe that in both options III and IV, more recent iterates (i.e., iterates $\mathbf{x}_{s,t}$ with larger indices t) will have larger contributions to \mathbf{x}^{s+1} . Theoretically, these two schemes can be analyzed in a unified manner, as shown in the following proposition.

Proposition 2. *In Algorithm 1, if we choose $\eta < \min\{b/12, 1\} / (\Gamma \bar{L})$ and m sufficiently large, then with either option III or IV, we have*

$$\mathbb{E}[f(\mathbf{x}^s) - f(\mathbf{x}^*)] \leq \bar{\rho}^s (f(\mathbf{x}^0) - f(\mathbf{x}^*)), \text{ where} \quad (31)$$

$$\bar{\rho} \triangleq \frac{b}{\gamma\bar{\mu}c'\eta(b - 4\eta\Gamma\bar{L}/(1 - \eta\gamma\bar{\mu}))} + \frac{4\eta\Gamma\bar{L}}{b - 4\eta\Gamma\bar{L}/(1 - \eta\gamma\bar{\mu})} \left(1 + \frac{1}{c'}\right) < 1 \quad (32)$$

and $c' \triangleq c/(1 - \eta\gamma\bar{\mu})^m$.

Remark 8. In the literature (Gower et al., 2016; Moritz et al., 2016), usually option I or II (in Algorithm 1) is analyzed theoretically to prove that the stochastic L-BFGS algorithms therein converge linearly. However, for faster convergence in practice, \mathbf{x}^{s+1} is chosen to be the last inner iterate $\mathbf{x}_{s,m}$. However, the latter is not amenable to linear convergence analysis. This *gap between theory and practice* is filled in by our geometric sampling or averaging scheme, i.e., option III or IV. Specifically, as shown in

Figure 2, our scheme not only yields linear convergence in theory, but also performs as well as the “last inner iterate” scheme in practice.

6.2 SUBSAMPLED GRADIENT STABILIZATION

In Algorithm 1, at the beginning of each outer iteration s , we compute a full gradient \mathbf{g}_s to stabilize the subsequent (inner) iterations. Inspired by Harikandeh et al. (2015), we propose a strategy that only computes a subsampled gradient $\tilde{\mathbf{g}}_s$ at the start of each outer iteration s . Specifically, we uniformly sample a subset $\tilde{\mathcal{B}}_s$ of $[n]$ with size \tilde{b}_s *without replacement* and then form $\tilde{\mathbf{g}}_s \triangleq (1/\tilde{b}_s) \sum_{i \in \tilde{\mathcal{B}}_s} \nabla f_i(\mathbf{x}^s)$. The size of $\tilde{\mathcal{B}}_s$, namely \tilde{b}_s , increases with the index s until it reaches n . By judiciously choosing \tilde{b}_s , we can show that the resulting algorithm still enjoys linear convergence with rate $\bar{\rho}$, when integrated with the geometric sampling/averaging scheme in Section 6.1. Before we formally state this result in Proposition 3, we first make an assumption in addition to Assumptions 1 and 2.

Assumption 5. *The inner iterates $\{\mathbf{x}_{s,t}\}_{s \geq 0, t \in [m]}$ generated by the modified algorithm in Section 6.2 are bounded almost surely, i.e., there exists $B < \infty$ such that for any $s \geq 0$ and $t \in [m]$, $\|\mathbf{x}_{s,t} - \mathbf{x}^*\| \leq B$.*

Proposition 3. *Let Assumptions 1, 2 and 5 hold. For any $\xi > 0$ and $S \in \mathbb{N}$, and for any $s \in [S]$, if we choose $\tilde{b}_s \geq \tilde{b}_s \triangleq nS^2\alpha_s / (S^2\alpha_s + (n-1)\xi^2\bar{\rho}^{2s})$, where $\alpha_s \triangleq 1/n \sum_{i=1}^n \|\nabla f_i(\mathbf{x}^s)\|^2$, we have*

$$\mathbb{E}[f(\mathbf{x}^s) - f(\mathbf{x}^*)] \leq \bar{\rho}^s \left\{ (f(\mathbf{x}^0) - f(\mathbf{x}^*)) + \left(1 + \frac{1}{c'}\right) \frac{\xi b}{b - 4\Gamma\bar{L}\eta/(1 - \eta\gamma\bar{\mu})} \left(\kappa_{\text{H}}^{1/2}B + \eta\Gamma\xi\right) \right\}.$$

Remark 9. Several remarks are in order. First, we remark that assumptions involving almost sure boundedness of iterates (e.g., Assumption 5) have appeared many times in the literature of stochastic optimization (Borkar, 2008; Harikandeh et al., 2015; Hu et al., 2009) and are always observed to hold in our experiments. Second, under this assumption, we can show that $\{\alpha_s\}_{s \geq 0}$ are bounded almost surely using Lipschitz continuity of ∇f_i , for any $i \in [n]$ in Assumption 2. Consequently, there exists $B' < \infty$ such that $\alpha_s \leq B'$ for any $s \geq 0$ and hence

$$\tilde{b}_s \leq \frac{nS^2B'}{S^2B' + (n-1)\xi^2\bar{\rho}^{2s}} \quad (33)$$

$$\leq \frac{nS^2B'}{(n-1)\xi^2} (\bar{\rho}^{-2})^s. \quad (34)$$

As a sanity check, we observe that (33) increases to n as $s \rightarrow \infty$. By further upper bounding (33) by (34), we obtain a practical rule to select \tilde{b}_s . Namely, it suffices to choose $\tilde{b}_s = \min\{\zeta v^s, n\}$, for some constants $\zeta > 0$

and $v > 1$. As shown in Section 7, this rule works well in practice. Third, for any $\epsilon > 0$, we can choose $S \in \mathbb{N}$ such that our algorithm achieves ϵ -suboptimality, i.e., $\mathbb{E}[f(\mathbf{x}^S) - f(\mathbf{x}^*)] < \epsilon$.

6.3 SMALL APPROXIMATE HESSIANS

In addition to high dimensionality and large size, *sparsity* is also a typical attribute for modern data, i.e., many feature vectors only have a few nonzero entries. For ERM problems (2) (with Tikhonov regularization), this implies that the Hessian of f in (1),

$$\nabla^2 f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ell''(\mathbf{a}_i^T \mathbf{x}, b_i) (\mathbf{a}_i \mathbf{a}_i^T) + \lambda \mathbf{I}, \quad (35)$$

tends to be sparse. Based on this observation, we propose a strategy that aims to approximate $\nabla^2 f(\mathbf{x})$ by several smaller Hessian matrices and update them efficiently. For sparse data, collecting curvature information via smaller dense Hessians can be more effective than directly manipulating the high-dimensional sparse Hessian. As a result, the algorithm converges faster in practice (see Figure 4).

Before describing our strategy, we first introduce some notations. We partition $[n]$ into K groups, and denote the set of partitions as $\mathcal{P} \triangleq \{\mathcal{P}_1, \dots, \mathcal{P}_K\}$. For any $i \in [K]$, we define $\mathcal{S}_i \triangleq \cup_{j \in \mathcal{P}_i} \text{supp}(\mathbf{a}_j)$, where $\text{supp}(\mathbf{a}_j)$ denotes the support of the vector \mathbf{a}_j . We define $d_i \triangleq |\mathcal{S}_i|$ and denote the elements in \mathcal{S}_i as $\{s_{i,1}, \dots, s_{i,d_i}\}$. We also define $F_i = \sum_{j \in \mathcal{P}_i} f_j$ so that $f = \frac{1}{n} \sum_{i=1}^K F_i$. We define a *projection matrix* $\mathbf{U}_i \in \mathbb{R}^{d_i \times d}$ such that for any $p \in [d_i]$ and $q \in [d]$, $u_{pq} = 1$ if $q = s_{i,p}$ and 0 otherwise. Accordingly, for any $l \in \mathcal{P}_i$, define a function $\phi_{i,l} : \mathbb{R}^{d_i} \rightarrow \mathbb{R}$ such that $f_l \triangleq \phi_{i,l} \circ \mathbf{U}_i$. Note that $\phi_{i,l}$ is uniquely defined by the definition \mathbf{U}_i . Also define $\phi_i \triangleq \sum_{l \in \mathcal{P}_i} \phi_{i,l}$. Therefore,

$$\nabla^2 f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^K \mathbf{U}_i^T \nabla^2 \phi_i(\mathbf{U}_i \mathbf{x}) \mathbf{U}_i, \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (36)$$

We now describe our strategy. In Algorithm 1, for any $i \in [K]$ and any $j \in \{r-M'+1, \dots, r\}$, define correction pairs $\mathbf{s}_{j,i} \triangleq \mathbf{U}_i(\bar{\mathbf{x}}_j - \bar{\mathbf{x}}_{j-1})$ and $\mathbf{y}_{j,i} \triangleq \sum_{l \in \mathcal{T}_{r,i}} \nabla^2 \phi_{i,l}(\mathbf{U}_i \bar{\mathbf{x}}_j) \mathbf{s}_{j,i}$, where $\mathcal{T}_{r,i}$ with size b_{H}/K is uniformly sampled from \mathcal{P}_i . Accordingly, define $\mathbf{S}_{r,i} \triangleq [\mathbf{s}_{r-M'+1,i}, \dots, \mathbf{s}_{r,i}]$, $\mathbf{Y}_{r,i} \triangleq [\mathbf{y}_{r-M'+1,i}, \dots, \mathbf{y}_{r,i}]$. Instead of storing \mathcal{H}_r , we only store matrices $\{(\mathbf{S}_{r,i}, \mathbf{Y}_{r,i})\}_{i=1}^K$. To reconstruct approximation $\mathbf{B}_{r,i}$ for each $\nabla^2 \phi_i$ at $\mathbf{U}_i \bar{\mathbf{x}}_r$, as usual, we apply M' BFGS updates (S-51) (using the correction pairs stored in $\mathbf{S}_{r,i}$ and $\mathbf{Y}_{r,i}$) to $\mathbf{B}_{r,i}^{(0)} \triangleq \delta_{r,i} \mathbf{I}$, where $\delta_{r,i} \triangleq \|\mathbf{y}_{r,i}\|^2 / \mathbf{s}_{r,i}^T \mathbf{y}_{r,i}$. This procedure can be implemented *efficiently* via a *compact representation* (Nocedal and Wright, 2006), i.e.,

$$\mathbf{B}_{r,i} = \delta_{r,i} \mathbf{I} - \mathbf{W}_{r,i} \mathbf{M}_{r,i}^{-1} \mathbf{W}_{r,i}^T, \quad (37)$$

$$\mathbf{M}_{r,i} \triangleq \begin{bmatrix} \delta_{r,i} \mathbf{S}_{r,i}^T \mathbf{S}_{r,i} & \mathbf{L}_{r,i} \\ \mathbf{L}_{r,i}^T & -\mathbf{D}_{r,i} \end{bmatrix}, \quad (38)$$

where $\mathbf{W}_{r,i} \triangleq [\delta_r^{(i)} \mathbf{S}_r^{(i)}, \mathbf{Y}_r^{(i)}]$ and $\mathbf{L}_{r,i}$ and $\mathbf{D}_{r,i}$ are the lower triangular matrix (excluding diagonal) and diagonal matrix of $\mathbf{S}_{r,i}^T \mathbf{Y}_{r,i}$ respectively. Analogous to (36), we define the approximation of $\nabla^2 f$ at $\mathbf{U}_i \bar{\mathbf{x}}_r$, denoted as \mathbf{B}_r , as

$$\mathbf{B}_r \triangleq \frac{1}{n} \sum_{i=1}^K \mathbf{U}_i^T \mathbf{B}_{r,i} \mathbf{U}_i. \quad (39)$$

We remark that the strong convexity of each function f_i (see Assumption 2), together with the full-row-rank property of \mathbf{U}_i , ensures $\nabla^2 \phi_i \succ 0$ on \mathbb{R}^{d_i} . This implies *positive curvature* $\mathbf{s}_{r,i}^T \mathbf{y}_{r,i} > 0$ and hence the positive definiteness of $\mathbf{B}_{r,i}$, for any $r \geq 0$ and $i \in [K]$. As a result, $\mathbf{B}_r \succ 0$ on \mathbb{R}^d . This suggests the usage of the *conjugate gradient* (CG) method to compute the search direction at time (s, t) , namely $\mathbf{p}_{s,t} \triangleq -\mathbf{B}_r^{-1} \mathbf{v}_{s,t}$, via solving the positive definite system $\mathbf{B}_r \mathbf{p}_{s,t} = -\mathbf{v}_{s,t}$. In particular, for any $\mathbf{z} \in \mathbb{R}^d$, $\mathbf{B}_r \mathbf{z}$ and $\mathbf{z}^T \mathbf{B}_r \mathbf{z}$ in CG can be computed very efficiently using (37) and (39). For example, $\mathbf{z}^T \mathbf{B}_r \mathbf{z}$ equals

$$\frac{1}{n} \sum_{i=1}^K \delta_{r,i} \|\mathbf{z}_i\|^2 - (\mathbf{W}_{r,i}^T \mathbf{z}_i)^T \mathbf{M}_{r,i}^{-1} (\mathbf{W}_{r,i}^T \mathbf{z}_i), \quad (40)$$

where $\mathbf{z}_i \triangleq \mathbf{U}_i \mathbf{z}$. We observe that the total computational cost in (40) is $O(M'(M'^2 + d'))$, where $d' \triangleq \sum_{i=1}^K d_i$. For sparse data, usually $d' = O(d)$, so this cost is still linear in d . In addition, we can compute (40) *in parallel* across $i \in [K]$. (Intuitively, this amounts to *collecting curvature* from each function ϕ_i in parallel.) In this case, the computational time will be greatly reduced to $O(M'(M'^2 + \max_i d_i))$. Since typically $\max_i d_i \ll d$, the computational savings from *parallel curvature collection* can be significant.⁹

7 NUMERICAL EXPERIMENTS

We consider two ERM problems, including logistic regression (with Tikhonov regularization) and ridge regression. For logistic regression, $b_i \in \{-1, 1\}$ and

$$f_i(\mathbf{x}) \triangleq \log \left(1 + e^{-b_i (\mathbf{a}_i^T \mathbf{x})} \right) + \frac{\lambda}{2} \|\mathbf{x}\|^2, \quad \forall i \in [n]. \quad (41)$$

For ridge regression, $b_i \in \mathbb{R}$ and

$$f_i(\mathbf{x}) \triangleq (\mathbf{a}_i^T \mathbf{x} - b_i)^2 + \frac{\lambda}{2} \|\mathbf{x}\|^2, \quad \forall i \in [n]. \quad (42)$$

Simple calculations show the smoothness parameters L_i for (41) and (42) are $\|\mathbf{a}_i\|^2 / 4 + \lambda$ and $2 \|\mathbf{a}_i\|^2 + \lambda$ respectively. In both (41) and (42), we choose $\lambda = 1/n$, following the convention in the literature (e.g., Konečný et al. (2016)).

We test logistic and ridge regression on `rcv1.binary` and `E2006-tfidf` datasets respectively (Chang and Lin, 2011). (In the sequel we abbreviate them as `rcv1` and `E2006`.) The size-dimension statistics (n, d) for these two

⁹The memory parameter M (note that $M' \leq M$) is usually set to a small constant, e.g., 5 or 10. Thus it has much less effects on the computational complexity compared to d or $\max_i d_i$.

datasets are (20242, 47236) and (16087, 15036) respectively. From (35), we define a sparsity estimate of $\nabla^2 f$ at any $\mathbf{x} \in \mathbb{R}^d$ as $\sigma \triangleq |\text{supp}(\mathbf{A}\mathbf{A}^T + \mathbf{I})| / d^2$, where $\mathbf{A} \triangleq [\mathbf{a}_1, \dots, \mathbf{a}_n]$ is the data matrix. For `rcv1` and `E2006`, σ equals 0.0154 and 0.0404 respectively. For both datasets, the norms of all feature vectors $\{\mathbf{a}_i\}_{i=1}^n$ have been normalized to unit. Since the the smoothness parameters L_i for both ERM problems only depend on $\|\mathbf{a}_i\|$ and λ , we have $L_i = L_j$ for any $i, j \in [n]$. Thus the nonuniform distribution P in Section 3 becomes uniform, and the advantage of nonuniform sampling of $\mathcal{B}_{s,t}$ cannot be observed.

To estimate the global optimum \mathbf{x}^* as ground truth, we use batch L-BFGS-B algorithm (Zhu et al., 1997). We randomly initialize \mathbf{x}^0 according to a scaled standard normal distribution.¹⁰ We use the number of data passes (i.e., number of data points accessed divided by n), rather than time, to measure the efficiency of algorithms. This has been a well-established convention in the literature since execution time highly depends on implementation of algorithms.

Finally we describe the parameter setting. We set the mini-batch size $b = \sqrt{n}$, Hessian update period $\Upsilon = 10$ and the memory parameter $M = 10$. We set $b_H = b\Upsilon$ so that the computation for \mathbf{y}_r can be amortized to each inner iteration. We set the number of inner iterations $m = n/b$, so that each outer iteration will access $2n$ data points. Lastly, we set $\eta = 1 \times 10^{-2}$. From Figure 1, we observe that when η is too large, e.g., $\eta = 0.1$, Algorithm 1 only converges sublinearly; whereas when η is too small, e.g., $\eta = 1 \times 10^{-3}$, Algorithm 1 converges linearly but slowly. This corresponds well to our theoretical analysis in Theorem 1, which indicates that when η falls below a threshold, ρ increases as η decreases. For both ERM problems, we find $\eta = 1 \times 10^{-2}$ yields fast linear convergence.

7.1 PERFORMANCE OF ACCELERATION STRATEGIES

We first examine the performance of Algorithm 1 with different schemes of choosing \mathbf{x}^{s+1} . We consider five schemes in total, including (a) uniform sampling (option I), (b) uniform averaging (option II), (c) geometric sampling (option III), (d) geometric averaging (option IV) and (e) last inner iterate (in Remark 8). For options III and IV, we set $\beta = 1/2$. From Figure 2, we observe that options III and IV perform as well as the “last inner iterate” scheme, on both ERM problems, and outperform the schemes based on uniform sampling/averaging significantly. For all the subsequent experiments, we use option IV to select \mathbf{x}^{s+1} .

We next compare the performance of Algorithm 1 with and without using the subsampled gradient stabilization strategy in Section 6.2. As suggested by Remark 9, we choose

¹⁰The performance of our algorithms were observed to be insensitive to the initialization of \mathbf{x}^0 .

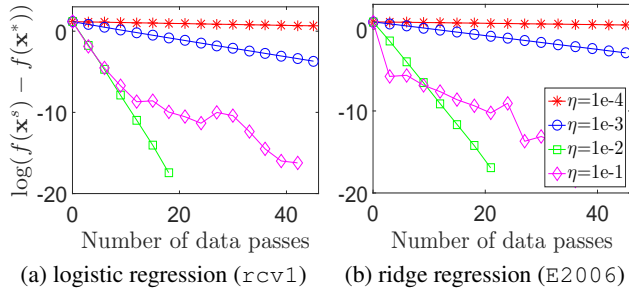


Figure 1: Log suboptimality versus number of passes through data of Algorithm 1 with different step sizes η .

$\tilde{b}_s = \min\{\zeta v^s, n\}$, where $\zeta = n/v^q$, $v = 3$ and $q = 8$. That is, the number of component gradients in $\tilde{\mathbf{g}}_s$ exponentially increases in the first $p = 8$ outer iterations and then remains at n . From Figure 3, we observe that this simple parameter selection method works well on both ERM problems, especially in the initial phase (when s is small). In addition, we also observe when s is large, both algorithms have almost the same (linear) convergence rates. This corroborates our analysis in Proposition 3.

We finally compare the performance of Algorithm 1 with and without using the low-dimensional approximate Hessian strategy in Section 6.3. We set the number of partitions $K = 5$ and partition $[n]$ evenly and randomly. From Figure 4, we observe that our strategy leads to improvements of convergence on both logistic and ridge regression problems, and the improvement on the latter is *very significant*. Moreover, our strategy *preserves the linear convergence* of Algorithm 1 on both problems. As discussed in Section 6.3, our strategy can be even more efficient in scenarios where parallel computational resources are available.

7.2 COMPARISON TO OTHER ALGORITHMS

We combine all of our acceleration strategies in Section 6 and compare the resulting algorithm with three benchmarking algorithms, including SVRG in Johnson and Zhang (2013) and two state-of-the-art stochastic L-BFGS algorithms in Moritz et al. (2016) and Gower et al. (2016). For the algorithm in Gower et al. (2016), we focus on its variant (b), since it consistently outperformed other variants in experiments. We tested all the three benchmarking algorithms on both ERM problems with different step sizes $\eta \in \{10^{-4}, 10^{-3}, \dots, 1\}$ and selected the best η for each algorithm. The outer iterate \mathbf{x}^{s+1} in all these algorithms are selected via “the last iterate” scheme. From Figure 5, we observe that on both ERM problems, our algorithm has faster convergence compared to all the benchmarking algorithms. The improvement of convergence is particularly significant on the ridge regression problem. This observation is consistent with our observations in Section 7.1.

Future work and open problems: see Section S-7.

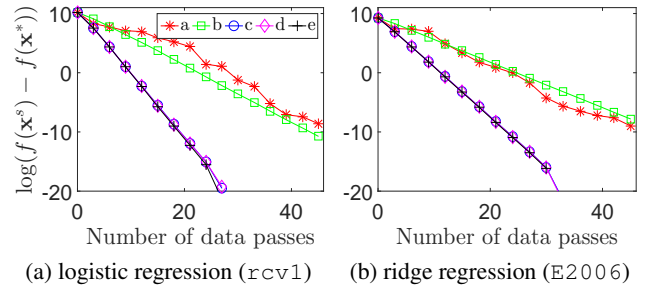


Figure 2: Comparison of Algorithm 1 with different selection schemes for \mathbf{x}^{s+1} .

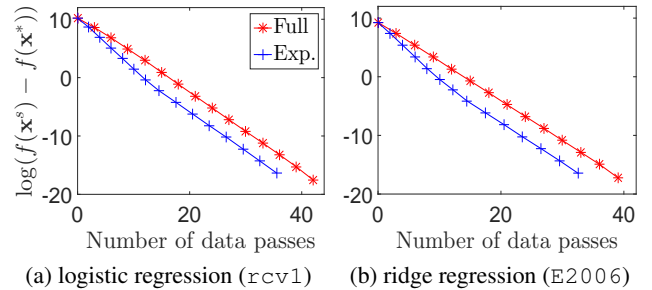


Figure 3: Comparison of Algorithm 1 without (Full) and with (Exp.) using the strategy in Section 6.2.

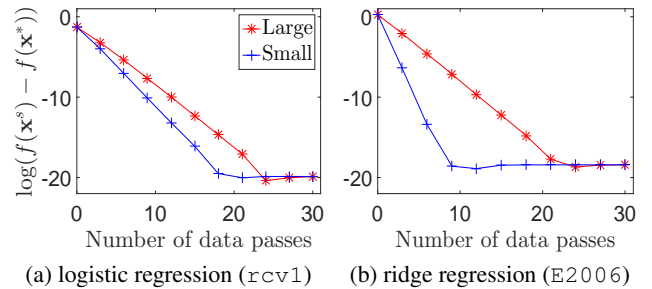


Figure 4: Comparison of Algorithm 1 without (Large) and with (Small) using the strategy in Section 6.3.

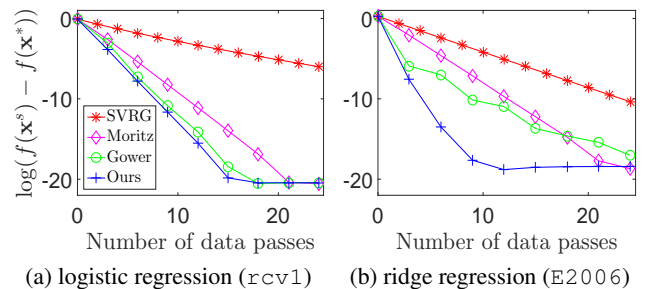


Figure 5: Comparison of our algorithm (Ours) with benchmarking algorithms (SVRG, Moritz and Gower).

References

- Z. Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. arXiv:1603.05953, 2016.
- A. Bordes, L. Bottou, and P. Gallinari. SGD-QN: Careful quasi-newton stochastic gradient descent. *J. Mach. Learn. Res.*, 10:1737–1754, Dec. 2009.
- V. S. Borkar. *Stochastic approximation: a dynamical systems viewpoint*. Cambridge, 2008.
- L. Bottou. Online algorithms and stochastic approximations. In *Online Learning and Neural Networks*. Cambridge University Press, 1998.
- L. Bottou and Y. LeCun. Large scale online learning. In *Proc. NIPS*, pages 1361–1368. Vancouver, BC, Canada, Dec. 2004.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- R. H. Byrd, S. L. Hansen, J. Nocedal, and Y. Singer. A stochastic quasi-newton method for large-scale optimization. *SIAM J. Optim.*, 26(2):1008–1031, 2016.
- C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.*, 2: 27:1–27:27, 2011.
- A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Proc. NIPS*, pages 1646–1654, Montréal, Québec, Canada, 2014.
- D. Goldfarb. A family of variable metric updates derived by variational means. *Math. Comput.*, 24(109):23–26, 1970.
- R. M. Gower, D. Goldfarb, and P. Richtárik. Stochastic block BFGS: squeezing more curvature out of data. In *Proc. ICML*, pages 1869–1878, New York City, NY, USA, June 2016.
- R. Harikandeh, M. O. Ahmed, A. Virani, M. Schmidt, J. Konečný, and S. Sallinen. Stopwasting my gradients: Practical svrg. In *Proc. NIPS*, pages 2251–2259. Montréal, Quebec, Canada, 2015.
- C. Hu, W. Pan, and J. T. Kwok. Accelerated gradient methods for stochastic optimization and online learning. In *Proc. NIPS*, pages 781–789. Vancouver, B.C., Canada, 2009.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Proc. NIPS*, pages 315–323, Lake Tahoe, Nevada, USA, 2013.
- J. Konečný and P. Richtárik. Semi-stochastic gradient descent methods. arXiv:1312.1666, 2013.
- J. Konečný, J. Liu, P. Richtárik, and M. Takáč. Mini-batch semi-stochastic gradient descent in the proximal setting. *IEEE J. Sel. Top. Signal Process.*, 10(2):242–255, 2016.
- H. Lin, J. Mairal, and Z. Harchaoui. A universal catalyst for first-order optimization. In *Proc. NIPS*, pages 3384–3392, 2015.
- D. C. Liu and J. Nocedal. On the limited memory bfgs method for large scale optimization. *Math. Program.*, 45(3), Dec. 1989.
- A. Mokhtari and A. Ribeiro. Global convergence of online limited memory BFGS. *J. Mach. Learn. Res.*, 16(1): 3151–3181, Jan 2015.
- P. Moritz, R. Nishihara, and M. I. Jordan. A linearly-convergent stochastic L-BFGS algorithm. In *Proc. AIS-TATS*, pages 249–258, Cadiz, Spain, May 2016.
- J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer, New York, 2nd edition, 2006.
- M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Math. Program.*, 162(1):83–112, 2017.
- N. N. Schraudolph, J. Yu, and S. Günter. A stochastic quasi-Newton method for online convex optimization. In *Proc. ATSTATS*, San Juan, Puerto Rico, March 2007.
- J. Sohl-Dickstein, B. Poole, and S. Ganguli. Fast large-scale optimization by unifying stochastic gradient and quasi-newton methods. In *Proc. ICML*, pages 604–612, Beijing, China, June 2014.
- D. Williams. *Probability with Martingales*. Cambridge University Press, 1991.
- L. Xiao and T. Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM J. Optim.*, 24(4):2057–2075, 2014.
- C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal. Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization. *ACM Trans. Math. Softw.*, 23(4), Dec 1997.